

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-244693

(43) 公開日 平成9年(1997)9月19日

(51) Int.Cl.⁶

G 1 0 L 5/04

識別記号

庁内整理番号

F I

G 1 0 L 5/04

技術表示箇所

D
F

審査請求 未請求 請求項の数 6 O L (全 7 頁)

(21) 出願番号 特願平8-49774

(22) 出願日 平成8年(1996)3月7日

(71) 出願人 000102728

エヌ・ティ・ティ・データ通信株式会社
東京都江東区豊洲三丁目3番3号

(72) 発明者 新村 貴彦

東京都江東区豊洲三丁目3番3号 エヌ・
ティ・ティ・データ通信株式会社内

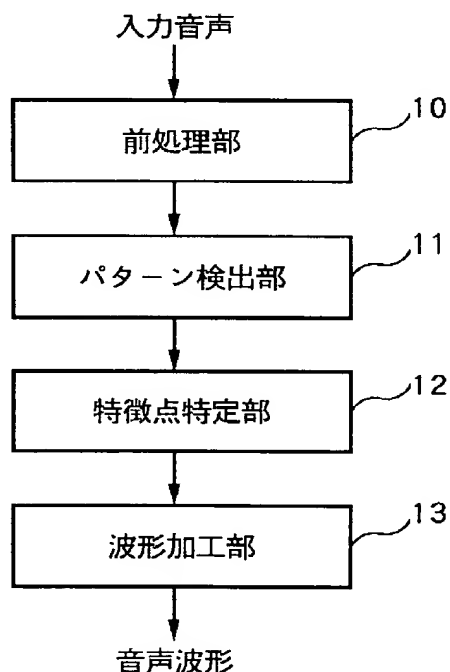
(74) 代理人 弁理士 鈴木 正剛

(54) 【発明の名称】 音声合成方法及び装置

(57) 【要約】

【課題】 音声合成装置において、音声モーフィング等の音声変換を行う際の合成音声の明瞭性を向上させる。

【解決手段】 声質の異なる二種以上の入力音声をそれぞれ音韻毎に区切り、対応する音韻毎に各入力音声を合成することで新たな音声を得る。その際、前処理部10で、文章を音韻単位に区切り、各音韻単位に属するピッチパターンをパターン検出部11で求める。特徴点特定部12では、ピッチパターンの特徴点、例えば極値のピッチ波形を入力音声毎に対応して特定する。波形加工部13では、対応するピッチ波形毎に合成処理を行って各音韻の合成処理を行う。



【特許請求の範囲】

【請求項1】 声質の異なる二種の入力音声をそれぞれ所定区間単位で区切り、個々の区間単位に対応するピッチ波形を組み合わせて合成音を生成する方法であって、個々の入力音声に属する前記区間単位内のピッチ波形のパターン（以下、ピッチパターン）の変化傾向を検出する過程と、

検出したピッチパターンの変化傾向の特徴点を区間単位毎に特定するとともに、各入力音声毎に、それぞれ特定された特徴点のピッチ波形及び特徴点間のピッチ波形に基づく新たなピッチパターンを生成する過程と、を含むことを特徴とする音声合成方法。

【請求項2】 前記生成された新たなピッチパターンに対応する元のピッチパターンに所定係数を乗じて新たなピッチ波形を生成する過程を含むことを特徴とする請求項1記載の音声合成方法。

【請求項3】 互いに声質の異なる第一及び第二の入力音声に対し、第一の入力音声から第二の入力音声へのモーフィング処理を行う方法であって、

前記第一及び第二の入力音声をそれぞれ所定区間単位で区切り、個々の区間単位内のピッチパターンの変化傾向を検出する過程と、

検出したピッチパターンの変化傾向の特徴点を区間単位毎に特定するとともに、各入力音声毎に、それぞれ特定された特徴点のピッチ波形及び特徴点間のピッチ波形に基づく新たなピッチパターンを生成する過程と、

生成された新たなピッチパターンに対応する元のピッチパターンに所定係数を乗じてモーフィング処理に用いるピッチ波形を生成する過程とを含む、

前記第二の入力音声に対する前記係数をモーフィングが進むにつれて大きくすることを特徴とする音声合成方法。

【請求項4】 前記特徴点が個々の区間単位におけるピッチパターンの極値であることを特徴とする請求項1ないし3のいずれかの項記載の音声合成方法。

【請求項5】 前記新たなピッチパターンにおける特徴点間のピッチ波形を、前記入力音声のピッチ波形の数をもとに補間することを特徴とする請求項1ないし4のいずれかの項記載の音声合成方法。

【請求項6】 声質の異なる二種の入力音声をそれぞれ所定区間単位で区切る前処理部と、

この前処理部で区切られた個々の区間単位内のピッチパターンの変化傾向を検出するパターン検出部と、

このパターン検出部で検出したピッチパターンの変化傾向の特徴点を区間単位毎に特定する特徴点特定部と、

この特徴点特定部により各入力音声毎に特定された特徴点のピッチ波形及び特徴点間のピッチ波形に基づく新たなピッチ波形を生成する波形加工部と、

を有し、この新たなピッチ波形を組み合わせて合成音を生成することを特徴とする音声合成装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、音声合成技術に関し、例えば、ある人物の音声を別人の音声へと変えていく音声モーフィング技術に関する。

【0002】

【従来の技術】音声合成技術は、駅構内でのアナウンスや機械による文章朗読等に広く用いられている。近年は、マルチメディア技術の台頭により、音声の表現力をより高めることが要求されてきており、例えば音声モーフィング処理等に代表される、声質変換という高度な技術が求められている。

【0003】音声モーフィング処理とは、画像処理技術において、ある人物の顔画像を徐々に別人の画像に変えていく過程を表現する画像モーフィング技術を音声に適用したもので、ある人物の声質を別人の声質へと変えていく過程を表現するものである。以下、音声モーフィングという場合、一の人間の声質を他の人間の声質に徐々に変えていくことをいうものとする。従来より、このような音声モーフィングは少なからず提案されており、例えば入力音声を音韻毎に区切って時間軸上で音声のスペクトrogram成分を線形に変化させるものが知られている。この技術において、読み上げ対象となる文を、最初的人物Aの音声により読み上げ、文の後半部は人物Bの音声に声質を変える場合の処理を以下に示す。

【0004】まず、人物A、Bのそれぞれに予め読み上げ対象の文を読み上げてもらい、そのピッチ波形を格納しておく。そして、読み上げ文の前半部は人物A、後半部は人物Bのピッチ波形をそのまま用い、中間部には、人物A、Bのピッチ波形を合成して得られる合成音声を用いて音声モーフィング処理を行う。

【0005】この処理の具体例を図8及び図9を用いて説明する。図8では「東京地方は多少雲が多いものの晴れ間も出ています」という文の読み上げの際に、音声モーフィング処理を行う場合の例を示すものである。この図に示されるように、文の前半部は人物Aの音声をそのまま用い、変換部分となる中間部、即ち「多い(おおい)」にあたる部分は合成音声により読み上げ、後半部を人物Bの音声をそのまま用いている。

【0006】図9に示されるように、上記音声の合成処理においては、まず人物の音韻(/o/, /o/, /i/)毎に音声を区切り、それぞれ実線矢印で示されるように、ピッチ波形の対応をとる。そして、対応するピッチ波形を足し合わせて、破線矢印で示される合成ピッチ波形を得る。次に、各音韻毎に、単純にピッチ波形の周波数の平均値を計算して点ピッチ周波数を求める。そして、人物A、Bの点ピッチ周波数の平均値を合成音のピッチ波形の周波数として設定する。さらに、合成したピッチ波形を、もとの語順に従って接続して合成音声を得る。

【0007】

【発明が解決しようとする課題】上述した従来の音声モーフィング処理においては、各音韻毎に合成音声のピッチ波形の周波数を算出している。この周波数は、ピッチ波形の間隔（時間長）の逆数となる。具体的には、図10（a）、（b）に示されるように、人物A、Bの音韻のピッチ波形毎に、図中の丸印で示される周波数をそれぞれ求め、さらに、図中の点線で示されるように、それぞれ点ピッチ周波数を算出している。合成音声のピッチ波形毎の周波数は、図10（c）に示されるように、人物A、Bの点ピッチ周波数の平均値を用い、各音韻について固定値としている。しかし、従来の音声モーフィング処理では、音韻をつなげて文章読み上げを行った場合に、その明瞭性が損なわれてしまうという難点が生じていた。

【0008】本発明の課題は、例えば音声モーフィング等の声質変換を行う際の合成音声の明瞭性を向上させる技術を提供することにある。

【0009】

【課題を解決するための手段】本発明者らによる検証の結果、従来の方法で明瞭性が劣るのは、音韻毎の点ピッチ周波数が固定値であり、点ピッチ周波数の変化が単調となる点に起因することが判明した。そこで、本発明では、音韻内での点ピッチ周波数を変化させる、改良された音声合成方法及び装置を創案した。

【0010】即ち、本発明の音声合成方法では、声質の異なる二種の入力音声をそれぞれ所定区間単位で区切り、個々の区間単位に対応するピッチ波形を組み合わせて合成音を生成する。その際、個々の入力音声に属する前記区間単位内のピッチパターンの変化傾向を検出する過程と、検出したピッチパターンの変化傾向の特徴点を区間単位毎に特定するとともに、各入力音声毎に、それぞれ特定された特徴点のピッチ波形及び特徴点間のピッチ波形に基づく新たなピッチパターンを生成する過程と、を含むことを特徴とする。

【0011】上記方法の後続処理としては、前記生成された新たなピッチパターンに対応する元のピッチパターンに所定係数を乗じて新たなピッチ波形を生成する過程が挙げられる。

【0012】本発明の他の音声合成方法は、互いに声質の異なる第一及び第二の入力音声に対し、第一の入力音声から第二の入力音声へのモーフィング処理を行う方法であって、前記第一及び第二の入力音声をそれぞれ所定区間単位で区切り、個々の区間単位内のピッチパターンの変化傾向を検出する過程と、検出したピッチパターンの変化傾向の特徴点を区間単位毎に特定するとともに、各入力音声毎に、それぞれ特定された特徴点のピッチ波形及び特徴点間のピッチ波形に基づく新たなピッチパターンを生成する過程と、生成された新たなピッチパターンに対応する元のピッチパターンに所定係数を乗じてモーフィング処理に用いるピッチ波形を生成する過程とを

含み、前記第二の入力音声に対する前記係数をモーフィングが進むにつれて大きくすることを特徴とする。

【0013】上記各方法の好ましい態様としては、前記特徴点を個々の区間単位におけるピッチパターンの極値とする。また、前記新たなピッチパターンにおける特徴点間のピッチ波形を、前記入力音声のピッチ波形の数をもとに補間する。

【0014】このようにして生成される音声は、元の入力音声のピッチパターンを反映して変化するものとなり、合成音の明瞭性が従来よりも向上する。特に、特徴点として、個々の区間単位におけるピッチ周波数の極値を用いることで、もとの音声の特徴が一層反映された音声生成される。さらに、新たなピッチパターンにおける特徴点間のピッチ波形を補間することで、元の入力音声の特徴を残しつつ、新たなピッチ波形が生成可能となる。

【0015】本発明は、また、上記各方法を実施する上で好適となる音声合成装置をも提供する。この装置は、声質の異なる二種の入力音声をそれぞれ所定区間単位で区切る前処理部と、この前処理部で区切られた個々の区間単位内のピッチパターンの変化傾向を検出するパターン検出部と、このパターン検出部で検出したピッチパターンの変化傾向の特徴点を区間単位毎に特定する特徴点特定部と、この特徴点特定部により各入力音声毎に特定された特徴点のピッチ波形及び特徴点間のピッチ波形に基づく新たなピッチ波形を生成する波形加工部と、を有し、この新たなピッチ波形を組み合わせて合成音を生成することを特徴とする。

【0016】

【発明の実施の形態】以下、図面を参照して本発明の実施の形態を詳細に説明する。この実施形態では、従来例と同様、「東京地方は多少雲が多いものの晴れ間もでています」という文章を例として説明する。

【0017】予め、この文章を人物A、Bにそれぞれ読み上げてもらい、その音声を入力音声としてそれぞれ格納しておく。この文章の前半の読み上げには、人物Aの音声をそのまま用い、後半の読み上げには、人物Bの音声をそのまま用いる。一方、この文章の中間部である「多い」の読み上げには、人物A、Bの音声から生成される音声をを用い、人物Aから人物Bへの声質の変換を行う。その際、中間部の音声は、図1に示す構成の音声合成装置により合成した。

【0018】この音声合成装置は、入力音声の波形を所定区間単位、例えば音韻を区間単位として区切る前処理部10と、個々の区間単位（音韻）内のピッチパターンの変化傾向を検出するパターン検出部11、検出したピッチパターンの変化傾向の特徴点を区間単位毎に特定する特徴点特定部12と、各入力音声毎に特定された特徴点のピッチ波形及び特徴点間のピッチ波形に基づく新たなピッチパターンを生成するとともに、新たなピッチ

パターンに対応する元のピッチパターンに所定係数、例えばある倍率数値を乗じて新たなピッチ波形を生成する波形加工部13とを備えている。

【0019】この実施形態では、図2(a)、(c)に示される人物A、Bの音韻に基づいて同(b)に示される合成音の音韻を得るため、上記ピッチパターンの変化傾向を特定するための特徴点として、ピッチパターンの極値を用いた。また、合成されたピッチ波形の周波数として、もとの各ピッチ波形の周波数の平均値を用いた。

【0020】以下、本実施形態による処理を、図3～図6をも参照して説明する。図3は、音声合成装置の全体的な処理概要を示すフローチャートであり、まず、前処理部10において、人物A、Bそれぞれの音声波形のうち、「多い」に相当する部分を音韻毎に区切る(S101)。音韻毎に区切られた音声波形を図4に示す。図4(a)は人物Aの音声波形、同(b)は人物Bの音声波形である。次に、パターン検出部11において、各音声波形から、ピッチ波形毎の周波数を求める(S102)。これはピッチ波形の間隔から容易に算出することができる。

【0021】その後、特徴点特定部12において、音韻の中で周波数が極値をとるピッチ波形を、人物A、Bの各音声波形毎に対応をとって特定する(S103、S104)。図5はこの様子を示すもので、(a)、(b)は人物Aの音韻のピッチパターン及びそのピッチ波形、(d)、(c)は人物Bの音韻のピッチパターン及びそのピッチ波形である。この例では人物Aのピッチ波形1と人物Bのピッチ波形1、人物Aのピッチ波形3と人物Bのピッチ波形5、及び人物Aのピッチ波形7と人物Bのピッチ波形9とがそれぞれ対応している。なお、他のピッチ波形に関しては、後述するように、もとの時系列の順に、適宜各ピッチ波形同士を対応させる(S105)。

【0022】次に、対応するピッチ波形を加算し、もとのピッチ波形の周波数の平均値を求めて合成ピッチ波形の周波数を決定する(S106)。この実施形態においては、図6に示されるように、各ピッチ波形に窓関数をかけたうえで加算することで、合成音声のピッチ波形を得た。また、振幅を正規化するため、振幅幅には係数である“0.5”を乗算した。

【0023】波形加工部13では、公知のピッチ波形重畳法を用いて、上述のようにして得られた合成ピッチ波形を、その周波数の間隔で並べ、図7に示される合成音を得る。このようにして、1つ1つのピッチ波形に対して、それぞれ独立に周波数が与えられる。

【0024】以上のように、本実施形態の音声合成装置では、音韻を構成するピッチ波形毎にそれぞれ周波数が与えられているので、音韻内での周波数が従来のように一定となってしまうことはない。特に、この音声合成装置においては、周波数が極値をとるピッチ波形同士を対

応させるようにしているので、各ピッチ波形における周波数が平均化されてしまうこともなく、周波数の高低差を損なうことなく合成音声を得ることができる。

【0025】なお、この音声合成装置においては、ピッチパターンが極値となるピッチ波形以外の各波形は、以下のように対応させることで補間した。まず、人物A、Bの音韻のピッチ波形毎に周波数を計算し、ピッチパターンの極値を求める。図5(a)、(d)に示されるように、この例においては、人物Aのピッチ波形は、1、3、7本目で極値を取り、人物Bのピッチ波形は、1、5、9本目で極値をとる。これら各ピッチ波形を、極値をとるピッチ波形によって、人物Aのピッチ波形は3本と4本に、人物Bの場合は5本と4本に、それぞれグループ分けする。

【0026】この際、人物Aの3本のピッチ波形と人物Bの5本のピッチ波形とを対応させるには、波形数の多いほうを分母、少ないほうを分子にして分数を決め、“1”から“5”までをかける。端数は切り捨て、“1”未満は“1”にする。これにより、人物A、Bによるピッチ波形の対応をとることができる。同様にして、人物Aの4本のピッチ波形と人物Bの4本のピッチ波形とを対応させる。以下、その具体的対応を示す。

【0027】まず、人物Bの1本目のピッチ波形は、 $1 \times 3 / 5 = 0.6$ なので人物Aの1本目のピッチ波形に対応させる。同様に、人物Bの2本目のピッチ波形は $2 \times 3 / 5 = 1.2$ なので人物Aの1本目、人物Bの3本目のピッチ波形は $3 \times 3 / 5 = 1.8$ なので人物Aの1本目、人物Bの4本目のピッチ波形は $4 \times 3 / 5 = 2.4$ なので人物Aの2本目、人物Bの5本目のピッチ波形は $5 \times 3 / 5 = 3.0$ なので人物Aの3本目、にそれぞれ対応させる。

【0028】次に、人物Bの6本目(Bの第二グループの一本目)のピッチ波形は $1 \times 4 / 4 = 1$ なので人物Aの4本目(Aの第二グループの一本目)に対応させる。人物Bの7本目のピッチ波形は $2 \times 4 / 4 = 2$ なので人物Aの5本目、人物Bの8本目のピッチ波形は $3 \times 4 / 4 = 3$ なので人物Aの6本目、人物Bの9本目のピッチ波形は $4 \times 4 / 4 = 4$ なので人物Aの7本目にそれぞれ対応させる。こうして、極値の対応を残したままで人物Aの音声と人物Bの音声との対応をとることができる。

【0029】なお、このような各波形の対応のさせ方については、特に限定がない。例えば、上記例では端数を切り捨てとしたが、端数を切り上げとしてもよいのはいうまでもない。また、この例では人物Aの声質から人物Bの声質への変形ステップが1ステップとなっているが、複数ステップにより声質を変形させてもよい。また、この実施形態では、もとのピッチ波形の点ピッチ周波数の単純平均値を合成音における周波数としているが、この周波数として、人物A、Bの点ピッチ周波数の加重平均値を用いてもよい。例えば人物Bの入力音声に

対する上記係数をモーフィングが進むにつれて大きくなるようにしてもよい。このようにして、合成音の前半は人物Aの周波数の影響を強くし、後半は人物Bの周波数の影響を強くすることで、人物Aから人物Bへの音声モーフィングをより一層滑らかに行うことも可能である。

【0030】

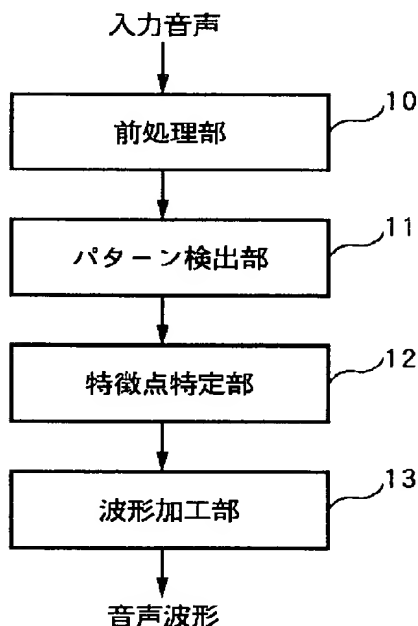
【発明の効果】以上の説明から明らかなように、本発明によれば、音韻内での点ピッチ周波数が一定となってしまうことはなく、従って、音韻の明瞭性が従来よりも向上する効果がある。特に、ピッチパターンの極値となるピッチ波形同士を対応させて波形加工を行うようにしたので、周波数の高いもの同士、低いもの同士でピッチ波形の合成がなされる。従って、周波数が高いピッチ波形と低いピッチ波形とが合成されることによる周波数の平均化現象が起こることもなくなり、音声の明瞭性が一層向上する。

【0031】また、音韻の周波数変化を利用して音声合成を行っているので、顔の動画像と組み合わせたときに、話者の唇の動きと細かな同期を取りやすくなる効果もある。

【図面の簡単な説明】

【図1】本発明の一実施形態の音声合成装置の要部ブロック構成図。

【図1】



【図2】(a)、(c)は人物A、Bの音声についての音韻のピッチパターンを示すグラフ、(b)は合成音についての音韻のピッチパターンを示すグラフ。

【図3】本実施形態による音声合成装置の処理概要を示すフローチャート。

【図4】(a)は合成対象となる音韻についての人物Aの音声波形、(b)は人物Bの音声波形を示す説明図。

【図5】(a)、(b)は人物Aについての音韻のピッチパターンとピッチ波形、(c)、(d)は人物Bについての音韻のピッチパターンとピッチ波形を示す説明図。

【図6】合成音のピッチ波形を得るまでの説明図。

【図7】合成された音声波形の説明図。

【図8】音声モーフィング処理の概要説明図。

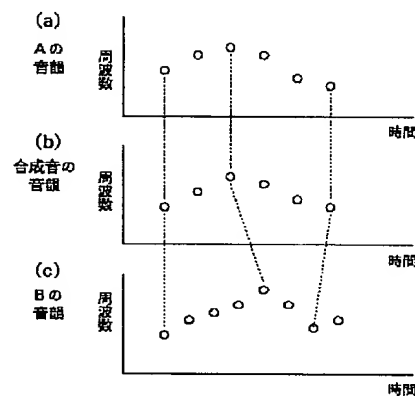
【図9】従来例における音声合成処理の手順説明図。

【図10】(a)、(b)は人物A、Bの音声についての音韻のピッチパターンを示すグラフ、(c)は合成音についての音韻のピッチパターンを示すグラフ。

【符号の説明】

- 10 前処理部
- 11 パターン検出部
- 12 特徴点特定部
- 13 波形加工部

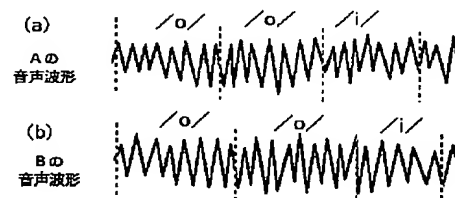
【図2】



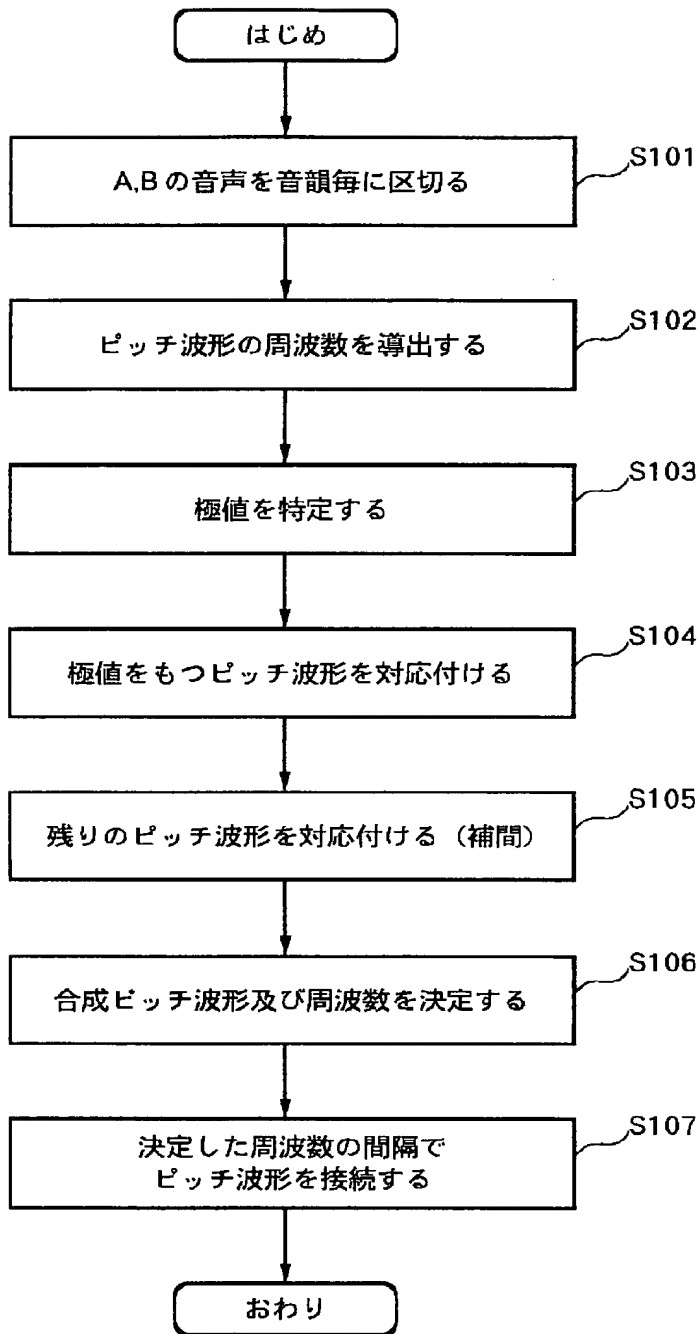
【図7】



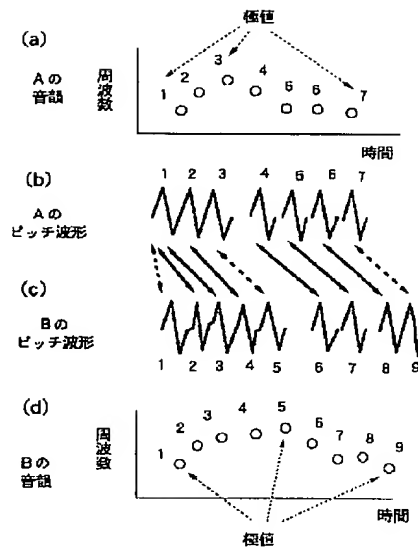
【図4】



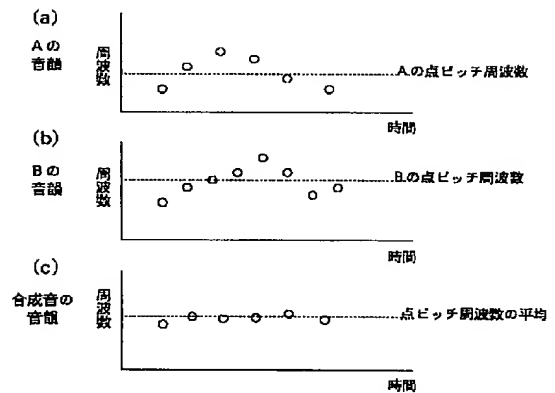
【図3】



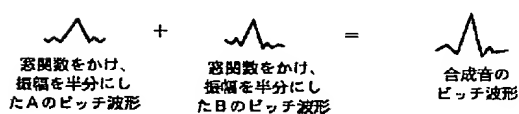
【図5】



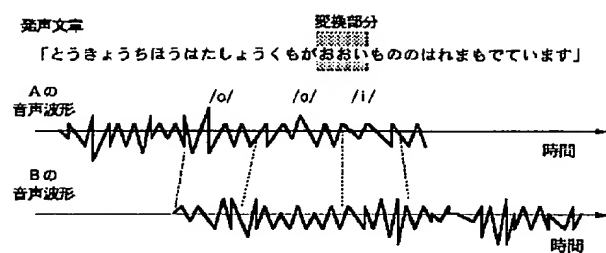
【図10】



【図6】



【図8】



【図9】

